



UTM
UNIVERSITI TEKNOLOGI MALAYSIA

FACULTY OF COMPUTING
SEMESTER II, SESSION 2024/2025

BACHELOR OF COMPUTER SCIENCE (BIOINFORMATICS)

SECB3203 PROGRAMMING FOR BIOINFORMATICS - SECTION 01

PROJECT
ALZHEIMER'S DISEASE PREDICTION

GROUP MEMBERS	MATRIC NO
TAN ZHAO HONG	A23CS0188
CHIN PEI WEN	A23CS0065
KOO XUAN	A23CS0300

LECTURER'S NAME : DR. SEAH CHOON SEN

SUBMISSION DATE : 03 DECEMBER 2025

Table of Contents

2.0 Data Collection and Pre-processing.....	3
2.1 Importing Dataset.....	4
- Python packages for Data Science.....	4
- Importing Dataset File.....	5
- Dataset Information.....	6
- Feature Removal.....	7
2.2 Data wrangling.....	8
- Target Variable Distribution.....	8
- Identifying and handling missing values.....	9
- Data Formatting.....	10
- Data Normalization.....	11
- Binning.....	13
- Indicator variables.....	15
- Train–Test Split and Z-Score Normalization.....	16

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB	AC	AD	AE	AF	AG	AH	AI	AJ	AK	AL	AM	AN	AO	AP	AQ	AR	AS	AT	AU	AV	AW	AX	AY	AZ	BA	BB	BC	BD	BE	BF	BG	BH	BI	BJ	BK	BL	BM	BN	BO	BP	BQ	BR	BS	BT	BU	BV	BW	BX	BY	BZ	CA	CB	CC	CD	CE	CF	CG	CH	CI	CJ	CK	CL	CM	CN	CO	CP	CQ	CR	CS	CT	CU	CV	CW	CX	CY	CZ	DA	DB	DC	DD	DE	DF	DG	DH	DI	DJ	DK	DL	DM	DN	DO	DP	DQ	DR	DS	DT	DU	DV	DW	DX	DY	DZ	EA	EB	EC	ED	EE	EF	EG	EH	EI	EJ	EK	EL	EM	EN	EO	EP	EQ	ER	ES	ET	EU	EV	EW	EX	EY	EZ	FA	FB	FC	FD	FE	FF	FG	FH	FI	FJ	FK	FL	FM	FN	FO	FP	FQ	FR	FS	FT	FU	FV	FW	FX	FY	FZ	GA	GB	GC	GD	GE	GF	GG	GH	GI	GJ	GK	GL	GM	GN	GO	GP	GQ	GR	GS	GT	GU	GV	GW	GX	GY	GZ	HA	HB	HC	HD	HE	HF	HG	HH	HI	HJ	HK	HL	HM	HN	HO	HP	HQ	HR	HS	HT	HU	HV	HW	HX	HY	HZ	IA	IB	IC	ID	IE	IF	IG	IH	II	IJ	IK	IL	IM	IN	IO	IP	IQ	IR	IS	IT	IU	IV	IW	IX	IY	IZ	JA	JB	JC	JD	JE	JF	JG	JH	JI	JJ	JK	JL	JM	JN	JO	JP	JQ	JR	JS	JT	JU	JV	JW	JX	JY	JZ	KA	KB	KC	KD	KE	KF	KG	KH	KI	KJ	KK	KL	KM	KN	KO	KP	KQ	KR	KS	KT	KU	KV	KW	KX	KY	KZ	LA	LB	LC	LD	LE	LF	LG	LH	LI	LJ	LK	LL	LM	LN	LO	LP	LQ	LR	LS	LT	LU	LV	LW	LX	LY	LZ	MA	MB	MC	MD	ME	MF	MG	MH	MI	MJ	MK	ML	MM	MN	MO	MP	MQ	MR	MS	MT	MU	MV	MW	MX	MY	MZ	NA	NB	NC	ND	NE	NF	NG	NH	NI	NJ	NK	NL	NM	NN	NO	NP	NQ	NR	NS	NT	NU	NV	NW	NX	NY	NZ	OA	OB	OC	OD	OE	OF	OG	OH	OI	OJ	OK	OL	OM	ON	OO	OP	OQ	OR	OS	OT	OU	OV	OW	OX	OY	OZ	PA	PB	PC	PD	PE	PF	PG	PH	PI	PJ	PK	PL	PM	PN	PO	PP	PQ	PR	PS	PT	PU	PV	PW	PX	PY	PZ	QA	QB	QC	QD	QE	QF	QG
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----

2.1 Importing Dataset

- Python packages for Data Science

```
4 import numpy as np
5 import pandas as pd
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8 from sklearn.model_selection import train_test_split, cross_val_score, GridSearchCV
9 from sklearn.preprocessing import StandardScaler, LabelEncoder
10 from sklearn.metrics import (accuracy_score, precision_score, recall_score,
11                               f1_score, confusion_matrix, classification_report,
12                               roc_auc_score, roc_curve)
13 # from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
14 from sklearn.linear_model import LogisticRegression
15 from sklearn.svm import SVC
16 from sklearn.tree import DecisionTreeClassifier
17 from sklearn.neighbors import KNeighborsClassifier
18 import warnings
19 warnings.filterwarnings('ignore')
```

- pandas: Data manipulation and analysis
- numpy: Numerical operations and array handling
- matplotlib: Data visualization
- seaborn: Statistical data visualization
- scikit-learn: Machine learning algorithms and preprocessing

- Importing Dataset File

```
37 # =====
38 # 1. DATA LOADING
39 # =====
40
41 def load_data(filepath):
42     """Load the Alzheimer's disease dataset"""
43     df = pd.read_csv(filepath)
44     print(f"Dataset shape: {df.shape}")
45     print(df)

```

```
367 if __name__ == "__main__":
368     FILEPATH = "alzheimers_disease_data.csv"
369     TARGET_COLUMN = 'Diagnosis'

```

```
=====
ALZHEIMER'S DISEASE PREDICTION - ML PIPELINE
=====

[1] Loading Data...
Dataset shape: (2149, 35)

```

	PatientID	Age	Gender	Ethnicity	EducationLevel	BMI	...	Disorientation	PersonalityChanges	DifficultyCompletingTasks	Forgetfulness	Diagnosis	DoctorInCharge
0	4751	73	0	0	2	22.927749	...	0	0	1	0	0	XXXConfid
1	4752	89	0	0	0	26.827681	...	0	0	0	1	0	XXXConfid
2	4753	73	0	3	1	17.795882	...	1	0	1	0	0	XXXConfid
3	4754	74	1	0	1	33.800817	...	0	0	0	0	0	XXXConfid
4	4755	89	0	0	0	20.716974	...	0	1	1	0	0	XXXConfid
...
2144	6895	61	0	0	1	39.121757	...	0	0	0	0	1	XXXConfid
2145	6896	75	0	0	2	17.857903	...	0	0	0	0	1	XXXConfid
2146	6897	77	0	0	1	15.476479	...	0	0	0	0	1	XXXConfid
2147	6898	78	1	3	1	15.299911	...	0	0	0	1	1	XXXConfid
2148	6899	72	0	0	2	33.289738	...	1	0	0	1	0	XXXConfid

```
[2149 rows x 35 columns]
```

These screenshots show the initial dataset loading process.

- Dataset Information

```
46 print(f"\nDataset info:")
47 print(df.info())
```

```
Dataset info:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2149 entries, 0 to 2148
Data columns (total 35 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   PatientID                            2149 non-null   int64
1   Age                                  2149 non-null   int64
2   Gender                              2149 non-null   int64
3   Ethnicity                            2149 non-null   int64
4   EducationLevel                       2149 non-null   int64
5   BMI                                  2149 non-null   float64
6   Smoking                              2149 non-null   int64
7   AlcoholConsumption                   2149 non-null   float64
8   PhysicalActivity                      2149 non-null   float64
9   DietQuality                          2149 non-null   float64
10  SleepQuality                         2149 non-null   float64
11  FamilyHistoryAlzheimers              2149 non-null   int64
12  CardiovascularDisease                2149 non-null   int64
13  Diabetes                             2149 non-null   int64
14  Depression                           2149 non-null   int64
15  HeadInjury                           2149 non-null   int64
16  Hypertension                         2149 non-null   int64
17  SystolicBP                           2149 non-null   int64
18  DiastolicBP                          2149 non-null   int64
19  CholesterolTotal                      2149 non-null   float64
20  CholesterolLDL                       2149 non-null   float64
21  CholesterolHDL                       2149 non-null   float64
22  CholesterolTriglycerides              2149 non-null   float64
23  MMSE                                 2149 non-null   float64
24  FunctionalAssessment                  2149 non-null   float64
25  MemoryComplaints                     2149 non-null   int64
26  BehavioralProblems                    2149 non-null   int64
27  ADL                                   2149 non-null   float64
28  Confusion                            2149 non-null   int64
29  Disorientation                        2149 non-null   int64
30  PersonalityChanges                   2149 non-null   int64
31  DifficultyCompletingTasks             2149 non-null   int64
32  Forgetfulness                         2149 non-null   int64
33  Diagnosis                             2149 non-null   int64
34  DoctorInCharge                       2149 non-null   object
dtypes: float64(12), int64(22), object(1)
memory usage: 587.7+ KB
None
```

These screenshots show the dataset information details.

- Feature Removal

```
25 # =====
26 # FEATURES TO REMOVE (Non-medical and weak evidence features)
27 # =====
28 FEATURES_TO_REMOVE = [
29     'PatientID',          # Administrative identifier
30     'DoctorInCharge',     # Administrative data
31     'Gender',             # Weak predictor
32     'Ethnicity',          # May introduce bias, weak medical relevance
33     'AlcoholConsumption', # Inconsistent evidence
34     'SleepQuality'        # Weak/inconsistent evidence
35 ]
36
```

```
51 # Remove non-medical and weak evidence features
52 print(f"\n{'='*80}")
53 print("REMOVING NON-MEDICAL AND WEAK EVIDENCE FEATURES")
54 print(f"{'='*80}")
55 features_removed = [col for col in FEATURES_TO_REMOVE if col in df.columns]
56 if features_removed:
57     print(f"\nRemoving features: {features_removed}")
58     df = df.drop(columns=features_removed)
59     print(f"New dataset shape: {df.shape}")
60 else:
61     print("\nNo features to remove (features not found in dataset)")
62
63 return df
```

```
=====
REMOVING NON-MEDICAL AND WEAK EVIDENCE FEATURES
=====
```

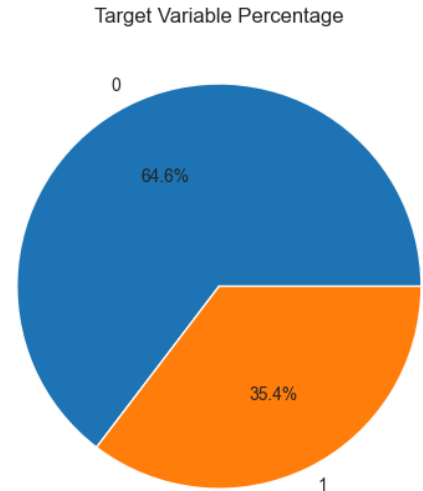
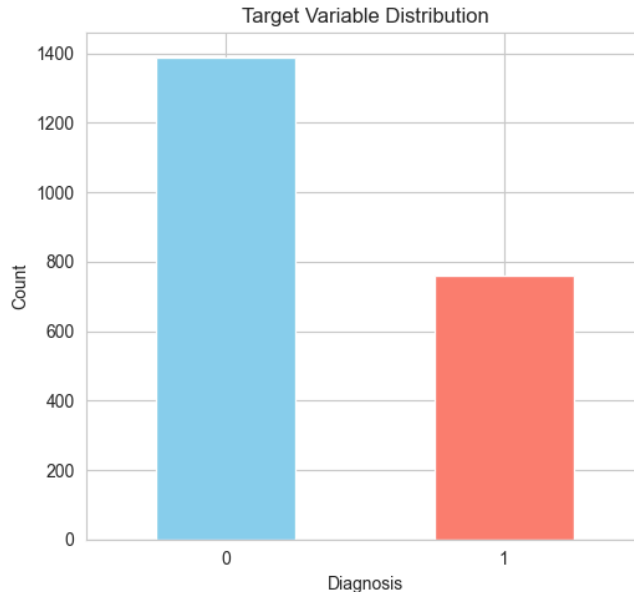
```
Removing features: ['PatientID', 'DoctorInCharge', 'Gender', 'Ethnicity', 'AlcoholConsumption', 'SleepQuality']
New dataset shape: (2149, 29)
```

These screenshots show the removal of non-medical and weak-evidence features such as patient identifiers and administrative attributes. These features were excluded to reduce noise, prevent bias, and ensure that the model focuses only on medically relevant variables that contribute meaningfully to Alzheimer's disease prediction.

2.2 Data wrangling

- Target Variable Distribution

```
110 def preprocess_data(df, target_column='Diagnosis'):
111     """Preprocess the data for machine learning"""
112
113     # Separate features and target
114     x = df.drop(columns=[target_column])
115     y = df[target_column]
116
117     # Encode target variable if it's categorical
118     if y.dtype == 'object':
119         print(f"\nEncoding target variable")
120         le_target = LabelEncoder()
121         y = le_target.fit_transform(y)
```



This figure illustrates the distribution of the target variable (Diagnosis) using both bar and pie charts. Visualizing the class distribution helps identify potential class imbalance and confirms that the problem is a binary classification task, which is suitable for the selected machine learning models.

- Identifying and handling missing values

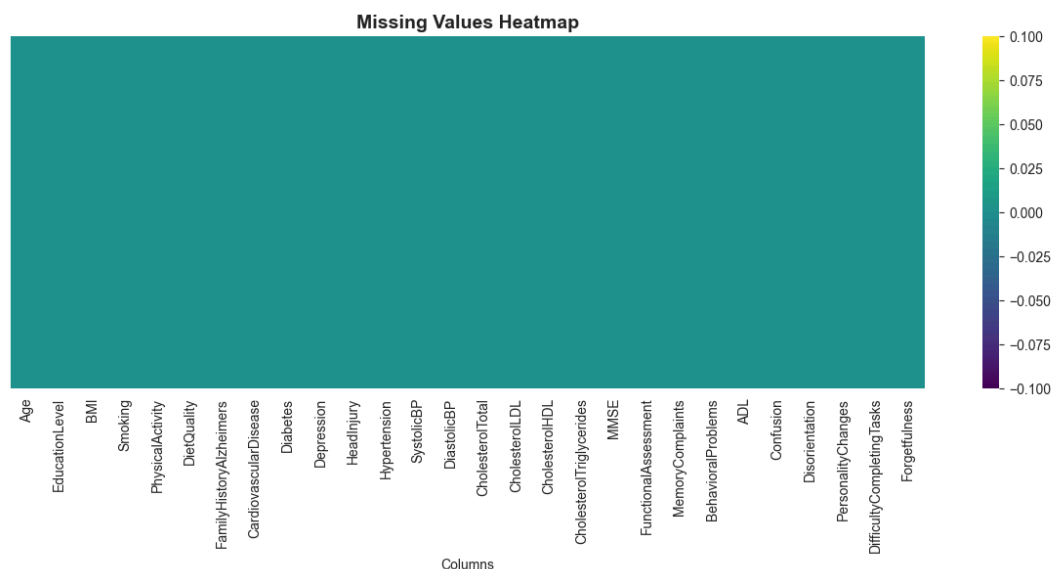
```

123 # VISUALIZE MISSING VALUES
124 print("\n--- Identifying and Handling Missing Values ---")
125 print("\n1. Missing Values Count:")
126 print(X.isnull().sum())
127
128 print("\n2. Visualizing Missing Values:")
129 plt.figure(figsize=(12, 6))
130 sns.heatmap(X.isnull(), cbar=True, cmap='viridis', yticklabels=False)
131 plt.title('Missing Values Heatmap', fontsize=14, fontweight='bold')
132 plt.xlabel('Columns')
133 plt.tight_layout()
134 plt.savefig('missing_values_heatmap.png', dpi=300)
135 plt.show()
136
137 # Handle missing values
138 if X.isnull().sum().sum() > 0:
139     print("\n3. Handling missing values...")
140     X = X.fillna(X.median())
141     print("Missing values after handling:")
142     print(X.isnull().sum())
143 else:
144     print("\n3. No missing values detected!")

```

2. Visualizing Missing Values:

3. No missing values detected!



These screenshots display the detection of missing values across all input features. Median imputation was applied to numerical attributes to handle missing data, as it is robust to outliers and preserves the overall distribution of the dataset. This step ensures data completeness before further processing.

- Data Formatting

```
146 # DATA FORMATTING
147 print("\n--- Data Formatting ---")
148 print("\n1. Current Data Types:")
149 print(x.dtypes)
150
151 print("\n2. Numerical Columns:")
152 numerical_cols = x.select_dtypes(include=[np.number]).columns
153 print(list(numerical_cols))
154
155 print("\n3. Categorical Columns:")
156 categorical_cols_orig = x.select_dtypes(include=['object']).columns
157 print(list(categorical_cols_orig))
```

```
--- Data Formatting ---
1. Current Data Types:
Age                int64
EducationLevel     int64
BMI                float64
Smoking            int64
PhysicalActivity   float64
DietQuality        float64
FamilyHistoryAlzheimers int64
CardiovascularDisease int64
Diabetes           int64
Depression         int64
HeadInjury         int64
Hypertension       int64
SystolicBP         int64
DiastolicBP        int64
CholesterolTotal   float64
CholesterolLDL     float64
CholesterolHDL     float64
CholesterolTriglycerides float64
MMSE               float64
FunctionalAssessment float64
MemoryComplaints   int64
BehavioralProblems int64
ADL                float64
Confusion          int64
Disorientation     int64
PersonalityChanges int64
DifficultyCompletingTasks int64
Forgetfulness      int64
dtype: object

2. Numerical Columns:
['Age', 'EducationLevel', 'BMI', 'Smoking', 'PhysicalActivity', 'DietQuality', 'FamilyHistoryAlzheimers', 'CardiovascularDisease', 'Diabetes', 'Depression', 'HeadInjury', 'Hypertension', 'SystolicBP', 'DiastolicBP', 'CholesterolTotal', 'CholesterolLDL', 'CholesterolHDL', 'CholesterolTriglycerides', 'MMSE', 'FunctionalAssessment', 'MemoryComplaints', 'BehavioralProblems', 'ADL', 'Confusion', 'Disorientation', 'PersonalityChanges', 'DifficultyCompletingTasks', 'Forgetfulness']

3. Categorical Columns:
[]
```

This step verifies and formats the dataset by inspecting data types and separating numerical and categorical features. Ensuring correct data types is essential before applying normalization, binning, and encoding, as machine learning algorithms require numerical and consistently formatted inputs.

- Data Normalization

```
159 # DATA NORMALIZATION
160 print("\n--- Data Normalization ---")
161
162 # Import MinMaxScaler
163 from sklearn.preprocessing import MinMaxScaler
164
165 print("\n1. Min-Max Scaling (Normalization)")
166 print("    Formula:  $x_{scaled} = (x - x_{min}) / (x_{max} - x_{min})$ ")
167 print("    Result: Values scaled to range [0, 1]")
168
169 if 'Age' in X.columns:
170     # Create a copy for demonstration
171     min_max_scaler = MinMaxScaler()
172     age_minmax = min_max_scaler.fit_transform(X[['Age']])
173
174     print("\n    Example: Age Column")
175     print(f"    Original range: [{X['Age'].min():.2f}, {(variable) age_minmax: ndarray[_AnyShape, dtype[Any]]}]")
176     print(f"    Scaled range: [{age_minmax.min():.2f}, {age_minmax.max():.2f}]")
177
178 # Show comparison
179 comparison_df = pd.DataFrame({
180     'Original_Age': X['Age'].head(10),
181     'MinMax_Scaled': age_minmax[:10].flatten()
182 })
183 print("\n    Sample comparison:")
184 print(comparison_df)
185
186 # Visualize
187 fig, axes = plt.subplots(1, 2, figsize=(14, 5))
188 axes[0].hist(X['Age'], bins=20, color='skyblue', edgecolor='black')
189 axes[0].set_title('Original Age Distribution', fontsize=12, fontweight='bold')
190 axes[0].set_xlabel('Age')
191 axes[0].set_ylabel('Frequency')
192
193 axes[1].hist(age_minmax, bins=20, color='salmon', edgecolor='black')
194 axes[1].set_title('Min-Max Scaled Age [0,1]', fontsize=12, fontweight='bold')
195 axes[1].set_xlabel('Scaled Age')
196 axes[1].set_ylabel('Frequency')
197
198 plt.tight_layout()
199 plt.savefig('minmax_scaling.png', dpi=300)
200 plt.show()
201
202 print("\n2. Z-Score Normalization (Standardization)")
203 print("    Formula:  $z = (x - \mu) / \sigma$ ")
204 print("    Note: This will be applied later using StandardScaler")
```

--- Data Normalization ---

1. Min-Max Scaling (Normalization)

Formula: $x_{\text{scaled}} = (x - x_{\text{min}}) / (x_{\text{max}} - x_{\text{min}})$

Result: Values scaled to range [0, 1]

Example: Age Column

Original range: [60.00, 90.00]

Scaled range: [0.00, 1.00]

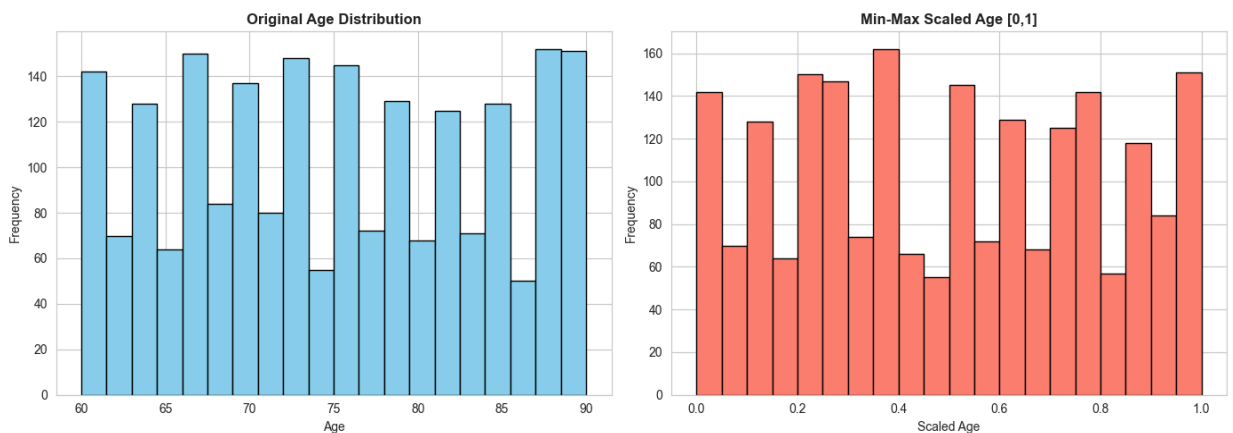
Sample comparison:

	Original_Age	MinMax_Scaled
0	73	0.433333
1	89	0.966667
2	73	0.433333
3	74	0.466667
4	89	0.966667
5	86	0.866667
6	68	0.266667
7	75	0.500000
8	72	0.400000
9	87	0.900000

2. Z-Score Normalization (Standardization)

Formula: $z = (x - \mu) / \sigma$

Note: This will be applied later using StandardScaler



This output demonstrates Min-Max normalization applied to the Age feature, scaling values into the range [0,1]. The comparison between original and scaled values highlights how normalization prevents features with larger numerical ranges from dominating the learning process of distance-based and gradient-based models.

- Binning

```
206 # BINNING
207 print("\n--- Binning ---")
208
209 if 'Age' in X.columns:
210     print("\n1. Age Binning:")
211     X['Age_binned'] = pd.cut(
212         X['Age'],
213         bins=[0, 60, 75, 100],
214         labels=['Young', 'Middle', 'Old']
215     )
216
217 # Show binning results
218 print("\n Sample of Age Binning:")
219 print(X[['Age', 'Age_binned']].head(10))
220
221 print("\n Age Group Distribution:")
222 print(X['Age_binned'].value_counts())
223
224 # Visualize
225 plt.figure(figsize=(10, 6))
226 X['Age_binned'].value_counts().sort_index().plot(kind='bar', color='steelblue', edgecolor='black')
227 plt.title('Age Group Distribution After Binning', fontsize=14, fontweight='bold')
228 plt.xlabel('Age Group')
229 plt.ylabel('Count')
230 plt.xticks(rotation=0)
231 plt.tight_layout()
232 plt.savefig('age_binning.png', dpi=300)
233 plt.show()
```

--- Binning ---

1. Age Binning:

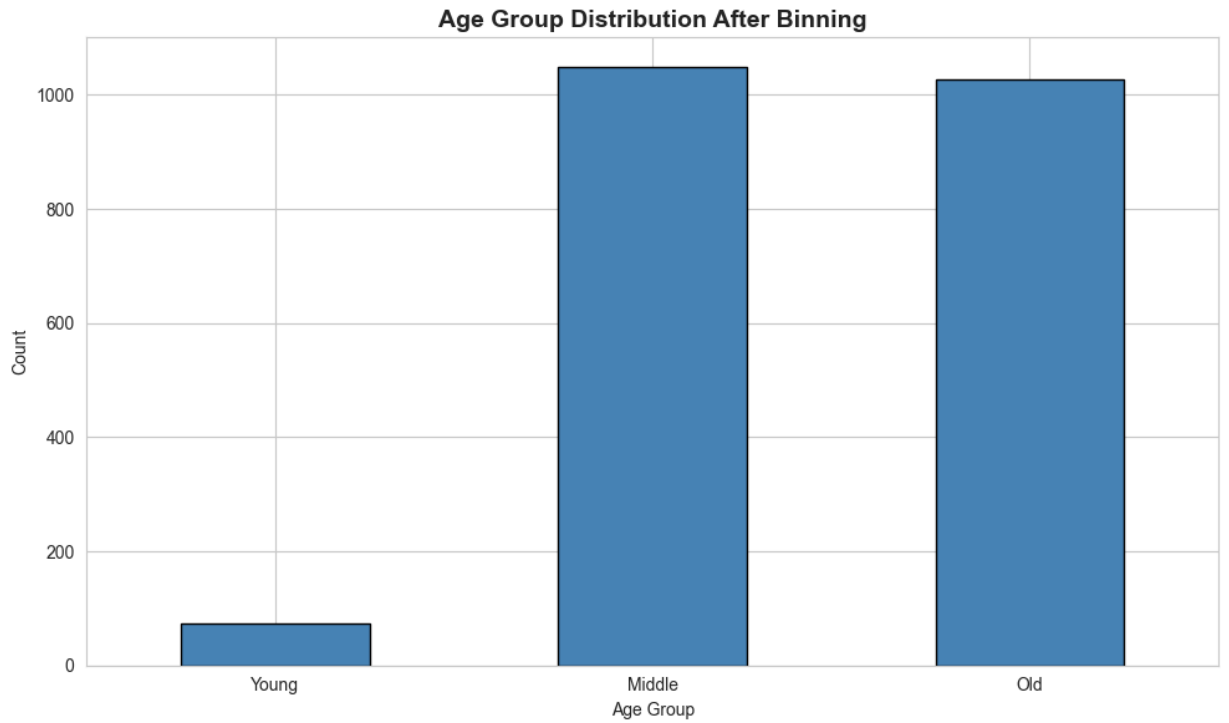
Sample of Age Binning:

	Age	Age_binned
0	73	Middle
1	89	Old
2	73	Middle
3	74	Middle
4	89	Old
5	86	Old
6	68	Middle
7	75	Middle
8	72	Middle
9	87	Old

Age Group Distribution:

Age_binned	
Middle	1048
Old	1027
Young	74

Name: count, dtype: int64



-

In this step, the continuous Age feature was transformed into categorical age groups using binning. This approach helps capture non-linear patterns related to age and Alzheimer's risk, improves interpretability, and allows the model to learn differences between broader age groups rather than relying only on exact numerical values.

- Indicator variables

```
235 # INDICATOR VARIABLES
236 print("\n--- Indicator Variables (One-Hot Encoding) ---")
237
238 categorical_cols = X.select_dtypes(include=['object', 'category']).columns
239
240 if len(categorical_cols) > 0:
241     print(f"\n1. Categorical columns to encode: {list(categorical_cols)}")
242
243     # Show before encoding
244     print("\n2. Before encoding (first 5 rows):")
245     print(X[categorical_cols].head())
246
247     # Apply one-hot encoding
248     X = pd.get_dummies(X, columns=categorical_cols, drop_first=True)
249
250     print(f"\n3. After one-hot encoding:")
251     print(f"    Original features: {len(categorical_cols)}")
252     print(f"    New indicator columns created: {len([col for col in X.columns if any(cat in col for cat in categorical_cols)])}")
253     print(f"    Total columns now: {X.shape[1]}")
```

```
--- Indicator Variables (One-Hot Encoding) ---

1. Categorical columns to encode: ['Age_binned']

2. Before encoding (first 5 rows):
   Age_binned
0      Middle
1        Old
2      Middle
3      Middle
4        Old

3. After one-hot encoding:
   Original features: 1
   New indicator columns created: 2
   Total columns now: 30
```

These screenshots show the conversion of categorical features created through binning into numerical indicator variables using one-hot encoding. This transformation is necessary because machine learning models require numerical inputs and cannot directly process categorical text labels.

- Train-Test Split and Z-Score Normalization

```
254
255     # SPLIT AND SCALE
256     # Split the data
257     X_train, X_test, y_train, y_test = train_test_split(
258         |     X, y, test_size=0.2, random_state=42, stratify=y
259     )
260
261     # Scale the features (Z-score normalization)
262     print("\n--- Final Scaling (Z-Score Normalization) ---")
263     scaler = StandardScaler()
264     X_train_scaled = scaler.fit_transform(X_train)
265     X_test_scaled = scaler.transform(X_test)
266
267     print(f"\nTraining set size: {X_train_scaled.shape}")
268     print(f"Test set size: {X_test_scaled.shape}")
269
270     print("\n" + "="*80)
271     print("DATA PREPROCESSING COMPLETE")
272     print("="*80)
273
274     return X_train_scaled, X_test_scaled, y_train, y_test, X.columns
275
```

```
--- Final Scaling (Z-Score Normalization) ---
```

```
Training set size: (1719, 30)
```

```
Test set size: (430, 30)
```

```
=====
DATA PREPROCESSING COMPLETE
=====
```

The dataset was split into training and testing sets using stratified sampling to preserve class distribution. Z-score normalization was then applied using StandardScaler to standardize features, ensuring that all variables contribute equally during model training.